

1. (10 points) The uniform distribution for a continuous variable x is

$$p(x; a, b) = \frac{1}{b-a}, \quad a \leq x \leq b.$$

Verify that this distribution is normalized (integrates to one), and find expressions for its mean and variance.

Awr-1 = Any distribution function is normalized if it integrates to 1 over its domain.

$$p(x; a, b) = \frac{1}{b-a} \quad \text{where } a \leq x \leq b$$

$$\int_a^b p(x; a, b) dx = \int_a^b \frac{1}{b-a} dx$$

$$= \frac{1}{b-a} [x]_a^b$$

$$= \frac{b-a}{b-a}$$

$$= 1$$

Hence this distribution is normalized.

Mean =

$$E(x) = \int_a^b x \cdot p(x) dx$$

$$= \int_a^b x \cdot \frac{1}{b-a} dx$$

$$= \frac{1}{b-a} \left[\frac{x^2}{2} \right]_a^b$$

$$= \frac{b^2 - a^2}{2(b-a)}$$

$$= \frac{(b-a)(b+a)}{2(b-a)}$$

$$= \frac{a+b}{2}$$

Variance =

$$E[(x - E[x^2])] = E[x^2] - (E[x])^2$$

$$\Rightarrow \int_a^b \frac{1}{b-a} x^2 - dx - \left(\frac{a+b}{2}\right)^2$$

[$E[x]$ = mean we found earlier]

$$\Rightarrow \frac{1}{b-a} \left[\frac{x^3}{3} \right]_a^b - \left(\frac{a+b}{2}\right)^2$$

$$= \frac{b^3 - a^3}{3(b-a)} - \left(\frac{a+b}{2}\right)^2$$

$$= \frac{(b-a)(b^2 + ab + a^2)}{3(b-a)} - \left(\frac{a^2 + b^2 + 2ab}{4}\right)$$

$$= \frac{b^2 + ab + a^2}{3} - \left(\frac{a^2 + b^2 + 2ab}{4}\right)$$

$$= \frac{4b^2 + 4ab + 4a^2 - 3a^2 - 3b^2 - 6ab}{12}$$

$$= \frac{a^2 + b^2 - 2ab}{12} = \frac{(a-b)^2}{12}$$

2. (10 points) Recall that the PMF of a Poisson random variable is

$$p(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!},$$

with one parameter λ . Given i.i.d. samples $x_1, \dots, x_n \sim \text{Pois}(\lambda)$, derive the MLE for λ .

Ans: 2 =

PMF of poisson random variable is =

$$p(x_i; \lambda) = \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

Taking log likelihood =

$$l(\lambda) = \ln \sum_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

$$= \sum_{i=1}^n \ln \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$= \sum_{i=1}^n (x_i \log \lambda - \lambda - \log x_i!)$$

$$= \log \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n \log x_i!$$

Equating derivative to 0 \Rightarrow

$$l'(r) = \frac{d}{dr} l(r) = 0$$

$$\Rightarrow \frac{d}{dr} \left\{ \log r \sum_{i=1}^n x_i - n r - \sum_{i=1}^n \log x_i! \right\} = 0$$

$$\Rightarrow \frac{1}{r} \sum_{i=1}^n x_i - n = 0$$

Therefore MLE = $\frac{\sum_{i=1}^n x_i}{n}$

3. (10 points) The function $\text{randn}(d, 1)$ generates a multivariate normal variable $\mathbf{x} \in \mathbb{R}^d$ with zero mean and covariance \mathbf{I} . Describe how to generate a random variable from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
Hint. If $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $A\mathbf{x} + \mathbf{b} \sim \mathcal{N}(A\boldsymbol{\mu} + \mathbf{b}, A\boldsymbol{\Sigma}A^\top)$.

$$\underline{\underline{\text{Ans. 3}}} = A \mathbf{x} \sim \mathcal{N}(0, I)$$

now as we have a matrix A and a normal variable \mathbf{x} . we will generate a random variable y such that $y = A\mathbf{x} + u$ and verify that $y \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

$$\text{then } A\mathbf{x} + b \sim \mathcal{N}(A\boldsymbol{\mu} + \mathbf{b}, A\boldsymbol{\Sigma}A^\top)$$

$$\text{Therefore for } \mathbf{x} \sim \mathcal{N}(0, I)$$

$$\text{Then } y = A\mathbf{x} + u \sim \mathcal{N}(A\mathbf{0} + u, AIA^\top)$$

$$\text{Hence } y \sim \mathcal{N}(\boldsymbol{\mu}, A\boldsymbol{\Sigma}A^\top) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

4. (20 points) Consider a data set in which each data sample i is associated with a weighting factor $r_i > 0$, and we instead try to minimize the weighted MSE function

$$\underset{\theta}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n r_i (y_i - \phi_i^\top \theta)^2.$$

Find an expression for the solution $\hat{\theta}$ that minimizes this loss function.

$$\underline{\text{Ans-4}} = \underset{\theta}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n r_i (y_i - \phi_i^\top \theta)^2$$

To minimize this weighted MSE function

we have; $\hat{\theta} = (\phi^\top \phi)^{-1} \phi^\top \psi$

where,

$$\psi = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

and

$$\phi = \begin{bmatrix} \phi_1^\top \\ \phi_2^\top \\ \vdots \\ \phi_n^\top \end{bmatrix}$$

Therefore we can write ϕ and ψ as

$$\psi' = \begin{bmatrix} \sqrt{r_1} y_1 \\ \sqrt{r_2} y_2 \\ \vdots \\ \sqrt{r_n} y_n \end{bmatrix} = P^\top \psi$$

$$\text{and } \phi' = \begin{bmatrix} \sqrt{r_1} \phi_1^\top \\ \sqrt{r_2} \phi_2^\top \\ \vdots \\ \sqrt{r_n} \phi_n^\top \end{bmatrix} = P^\top \phi$$

where $\rho = \begin{bmatrix} \sqrt{\sigma_1} & 0 & \cdots & 0 \\ 0 & \sqrt{\sigma_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{\sigma_n} \end{bmatrix}$

$$\begin{aligned}\hat{\Omega} &= (\phi'^T \phi')^{-1} \phi'^T \psi' \\ &= ((P^T \phi)^T (P^T \phi))^{-1} (P^T \phi)^T (P^T \psi) \\ &= (\phi^T P \cdot P^T \phi)^{-1} (\phi^T P \cdot P^T \psi)\end{aligned}$$

where,

$$P \cdot P^T = \begin{bmatrix} \sqrt{\sigma_1} & 0 & \cdots & 0 \\ 0 & \sqrt{\sigma_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{\sigma_n} \end{bmatrix} \quad \begin{bmatrix} \sqrt{\sigma_1} & 0 & \cdots & 0 \\ 0 & \sqrt{\sigma_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{\sigma_n} \end{bmatrix}$$

$$P \cdot P^T = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n \end{bmatrix}$$

5. (50 points) The 20 Newsgroups data set is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups. The data set can be downloaded here: <<http://qwone.com/~jason/20Newsgroups/>>. For simplicity, we will just focus on the “bag-of-words” representation of the documents given in the Matlab/Octave section. In this case, the input \mathbf{x}_i is a vector of word histogram of doc i , and the output y_i is the news group that it belongs to.

- (a) One effective classifier by Tom Mitchell is an instance of the naive Bayes model. First, the actual word count is ignored in the input data; we only consider whether a word j appears in doc i or not. Then each feature in \mathbf{x} can be viewed as a Bernoulli random variable. Furthermore, the naive Bayes assumption states that each of these Bernoulli random variables are conditionally independent given the label y , i.e., $p(\mathbf{x}|y) = \prod p(x_j|y)$. Each of the $p(x_j|y)$ can be easily estimated using the training data.

Ans. (a) =

Derivation for mathematical expression
for Bayes model.

$$\mathbf{x} \in \{0, 1\}^k$$

$$Y \sim \text{Bernoulli}(\phi)$$

$$x_k \sim \text{Bernoulli}(\theta_k | \gamma) \quad \forall k \in \{1, \dots, k\}$$

$$\text{Bern}(x - \theta) = P(x|\theta) = \theta^x (1-\theta)^{1-x}$$

$$\mathcal{D} = \{x_1, x_2, \dots, x_N\}$$

$$P(\theta|x) \propto P(\theta) P(x|\theta)$$

$$P(\theta | a, b) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(a, b)}$$

$$P(\theta | D) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \theta^{n_1} (1-\theta)^{n_2}$$

$$P(\theta | D) \propto \theta^{n_1 + \alpha - 1} (1-\theta)^{n_2 + \beta - 1}$$

$$P(\theta | D) = \frac{\theta^{n_1 + \alpha - 1} (1-\theta)^{n_2 + \beta - 1}}{B(\alpha + n_1 + b + n_2)}$$

$$\hat{\theta}_{\text{map}} = \frac{\alpha + n_1 - 1}{\alpha + \beta + n_1 + n_2 - 2} = \frac{n_1 + 1}{N + 2}$$

N = number of documents in class that contain feature n_1

n is the count of feature in class after getting $\hat{\theta}_{\text{map}}$ for all the features in vocabulary. we extract features for test data. we multiply the features with all the corresponding $\hat{\theta}_{\text{map}}$

for all values of k where k is the number of class.

The class score with maximum value is assigned as class of the data.

$$\hat{\theta}_{\text{map}} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N P(x^{(i)} | \theta) p(\theta)$$

- (b) To incorporate the word count, we can impose a rather different probabilistic model. Assume each doc is a huge multinomial random variable, with cardinality equal to the vocabulary size and the total number of draws is the length of that doc, given the label. In other words, $p(\mathbf{x}|y)$ is multinomial.

Ans-b = we know that,

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

MLE of $c \Rightarrow$

$$\begin{aligned} c_{\text{MAP}} &= \operatorname{argmax}_{c \in C} P(c|d) \\ &= \operatorname{argmax}_{c \in C} \frac{P(d|c)P(c)}{P(d)} \\ &= \operatorname{argmax}_{c \in C} P(d|c)P(c) \\ &= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c) \end{aligned}$$

$$\begin{aligned} c_{\text{NB}} &= \operatorname{argmax}_{c \in C} P(c_3) \prod_{x \in X} P(x|c) \\ &= \operatorname{argmax}_{c \in C} P(c_3) \prod_{i \in \text{position}} P(x_i | c_3) \end{aligned}$$

$$\hat{P}(C_3) = \frac{\text{Count Documents } (w_i, C_3)}{\sum_{\omega \in V} \text{count } (\omega, C_3)}$$

Parameter Smoothing =

$$\hat{P}(w_i | C_3) = \frac{\text{count } (w_i, C_3)}{\sum_{\omega \in V} \text{count } (\omega, C_3)}$$

Laplace Smoothing =

$$\hat{P}(w_i | C) = \frac{\text{count } (w_i, C_3) + 1}{\sum_{\omega \in V} (\text{count } (\omega, C_3) + 1)}$$

$$= \frac{\text{count } (w_i, C_3) + 1}{\sum_{\omega \in V} (\text{count } (\omega, C_3) + |V|)}$$